

Backpropagation

Deep Learning (DSE316/616)

Vinod K Kurmi
Assistant Professor, DSE

Indian Institute of Science Education and Research Bhopal

Aug 18, 2022



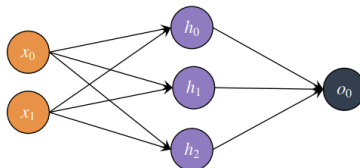
Disclaimer

- Much of the material and slides for this lecture were borrowed from
 - Bernhard Schölkopf's MLSS 2017 lecture,
 - Tommi Jaakkola's 6.867 class,
 - CMP784: Deep Learning Fall 2021 Erkut Erdem Hacettepe University
 - Fei-Fei Li, Andrej Karpathy and Justin Johnson's CS231n class
 - Hongsheng Li's ELEG5491 class

Previous Class: Recap

Input

[
[-20, 45],
[80, 0],
[4, 15],
[45, 60],
]

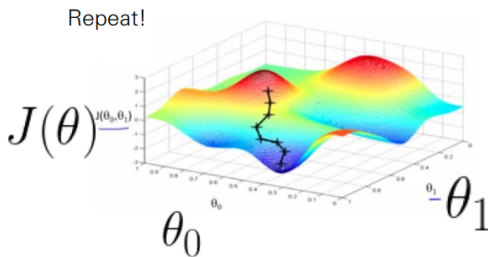


Predicted Actual

[[
0.05	1
0.02	0
0.96	1
0.35	1
]]

$$J(\theta) = \frac{1}{N} \sum_i \ell(\underbrace{f(\mathbf{x}^{(i)}; \theta)}_{\text{Predicted}}, \underbrace{y^{(i)}}_{\text{Actual}})$$

Previous Class: loss surface



Previous Class: SGD

- Initialize θ randomly
- For N Epochs
 - For each training batch $\{(x_0, y_0), (x_1, y_1), \dots (x_B, y_B)\}$:
 - Compute Loss
 - Gradient:
$$\frac{\partial \mathcal{J}(\theta)}{\partial \theta} = \frac{1}{B} \sum_i \frac{\partial \mathcal{J}_i(\theta)}{\partial \theta}$$
 - Update θ with update rule: $\theta = \theta - \eta \frac{\partial \mathcal{J}(\theta)}{\partial \theta}$

Advantages

- More accurate estimation of gradient
 - Smoother convergence
 - Allows for larger learning rates
- Minibatches lead to fast training!
 - Can parallelize computation

Previous Class: Training NN

- $a_L(x; \theta_1, \dots, \theta_L) = h_L(h_{L-1}(\dots h_1(x; \theta_1), \theta_{L-1}); \theta_L)$
 - x : input
 - θ_l : parameter of layer l
 - $a_l = h_l(x; \theta_l)$: (non)-linear function
- Given training corpus X, Y find optimal parameters

$$\theta^* = \underset{\theta}{\operatorname{argmin}} \sum_{(x,y) \in (X,Y)} \ell(y, a_L(x; \theta_1, \dots, \theta_L))$$

- To use any gradient descent based optimization

$$\theta^{(t+1)} = \theta^t - \eta_t \frac{\partial \mathcal{L}}{\partial \theta_t}$$

- we need the gradients $\frac{\partial \mathcal{L}}{\partial \theta_l}$; $l = 1, 2, \dots, L$
- How to compute the gradients for such a complicated function enclosing other functions, like $a_l(\dots)$

Training a neural network

- **Given:**
 - A network architecture (layout of neurons, their connectivity and activations)
 - A dataset of labeled examples $S = \{(x_i, y_i)\}$
- **The goal:** Learn the weights of the neural network
- **Remember:** For a fixed architecture, a neural network is a function parameterized by its weights
 - Prediction $y = NN(x; w)$

Training a neural network

- **Given:**
 - A network architecture (layout of neurons, their connectivity and activations)
 - A dataset of labeled examples $S = \{(x_i, y_i)\}$
- **The goal:** Learn the weights of the neural network
- **Remember:** For a fixed architecture, a neural network is a function parameterized by its weights
 - Prediction $y = NN(x; w)$

Training a neural network

- **Given:**
 - A network architecture (layout of neurons, their connectivity and activations)
 - A dataset of labeled examples $S = \{(x_i, y_i)\}$
- **The goal:** Learn the weights of the neural network
- **Remember:** For a fixed architecture, a neural network is a function parameterized by its weights
 - Prediction $y = NN(x; w)$

Training a neural network

- **Given:**
 - A network architecture (layout of neurons, their connectivity and activations)
 - A dataset of labeled examples $S = \{(x_i, y_i)\}$
- **The goal:** Learn the weights of the neural network
- **Remember:** For a fixed architecture, a neural network is a function parameterized by its weights
 - Prediction $y = NN(x; w)$

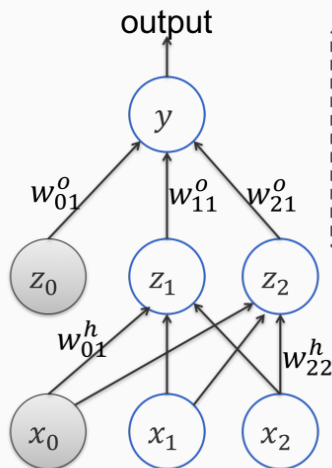
Training a neural network

- **Given:**
 - A network architecture (layout of neurons, their connectivity and activations)
 - A dataset of labeled examples $S = \{(x_i, y_i)\}$
- **The goal:** Learn the weights of the neural network
- **Remember:** For a fixed architecture, a neural network is a function parameterized by its weights
 - Prediction $y = NN(x; w)$

Training a neural network

- **Given:**
 - A network architecture (layout of neurons, their connectivity and activations)
 - A dataset of labeled examples $S = \{(x_i, y_i)\}$
- **The goal:** Learn the weights of the neural network
- **Remember:** For a fixed architecture, a neural network is a function parameterized by its weights
 - Prediction $y = NN(x; w)$

Back to our running example



Given an input \mathbf{x} , how is the output predicted

$$\text{output } y = w_{01}^o + w_{11}^o z_1 + w_{21}^o z_2$$

$$z_2 = \sigma(w_{02}^h + w_{12}^h x_1 + w_{22}^h x_2)$$

$$z_1 = \sigma(w_{01}^h + w_{11}^h x_1 + w_{21}^h x_2)$$

Suppose the true label for this example is a number y_i

We can write the *square loss* for this example as:

$$L = \frac{1}{2} (y - y_i)^2$$

Checkpoints

- If we have a neural network (structure, activations and weights), we can make a prediction for an input
- If we had the true label of the input, then we can define the loss for that example
- If we can take the derivative of the loss with respect to each of the weights, we can take a gradient step in SGD

Neural networks: the original linear classifier

- Linear score function: $\mathbf{f} = \mathbf{W}\mathbf{x}$
 - $x \in \mathcal{R}^D$
 - $W \in \mathcal{R}^{C \times D}$
- Two layer neural network: $f = W_2 \max(0, W_1 x)$
 - $x \in \mathcal{R}^D$
 - $W_1 \in \mathcal{R}^{H \times D}$
 - $W_2 \in \mathcal{R}^{C \times H}$
- In practice we will usually add a learnable bias at each layer as well
- Neural Network is a very broad term; these are more accurately called "fully-connected networks" or sometimes "multi-layer perceptrons" (MLP)

Neural networks: the original linear classifier

- Linear score function: $\mathbf{f} = \mathbf{W}\mathbf{x}$
 - $\mathbf{x} \in \mathcal{R}^D$
 - $\mathbf{W} \in \mathcal{R}^{C \times D}$
- Two layer neural network: $f = W_2 \max(0, W_1 x)$
 - $\mathbf{x} \in \mathcal{R}^D$
 - $\mathbf{W}_1 \in \mathcal{R}^{H \times D}$
 - $\mathbf{W}_2 \in \mathcal{R}^{C \times H}$
- In practice we will usually add a learnable bias at each layer as well
- Neural Network is a very broad term; these are more accurately called "fully-connected networks" or sometimes "multi-layer perceptrons" (MLP)

Neural networks: the original linear classifier

- Linear score function: $\mathbf{f} = \mathbf{W}\mathbf{x}$
 - $\mathbf{x} \in \mathcal{R}^D$
 - $\mathbf{W} \in \mathcal{R}^{C \times D}$
- Two layer neural network: $f = W_2 \max(0, W_1 x)$
 - $\mathbf{x} \in \mathcal{R}^D$
 - $\mathbf{W}_1 \in \mathcal{R}^{H \times D}$
 - $\mathbf{W}_2 \in \mathcal{R}^{C \times H}$
- In practice we will usually add a learnable bias at each layer as well
- Neural Network is a very broad term; these are more accurately called "fully-connected networks" or sometimes "multi-layer perceptrons" (MLP)

Neural networks: the original linear classifier

- Linear score function: $\mathbf{f} = \mathbf{W}\mathbf{x}$
 - $x \in \mathcal{R}^D$
 - $W \in \mathcal{R}^{C \times D}$
- Two layer neural network: $f = W_2 \max(0, W_1 x)$
 - $x \in \mathcal{R}^D$
 - $W_1 \in \mathcal{R}^{H \times D}$
 - $W_2 \in \mathcal{R}^{C \times H}$
- In practice we will usually add a learnable bias at each layer as well
- Neural Network is a very broad term; these are more accurately called "fully-connected networks" or sometimes "multi-layer perceptrons" (MLP)

Neural networks: the original linear classifier

- Linear score function: $\mathbf{f} = \mathbf{W}\mathbf{x}$
 - $\mathbf{x} \in \mathcal{R}^D$
 - $\mathbf{W} \in \mathcal{R}^{C \times D}$
- Two layer neural network: $f = W_2 \max(0, W_1 x)$
 - $\mathbf{x} \in \mathcal{R}^D$
 - $\mathbf{W}_1 \in \mathcal{R}^{H \times D}$
 - $\mathbf{W}_2 \in \mathcal{R}^{C \times H}$
- In practice we will usually add a learnable bias at each layer as well
- Neural Network is a very broad term; these are more accurately called "fully-connected networks" or sometimes "multi-layer perceptrons" (MLP)

Neural networks: the original linear classifier

- Linear score function: $\mathbf{f} = \mathbf{W}\mathbf{x}$
 - $x \in \mathcal{R}^D$
 - $W \in \mathcal{R}^{C \times D}$
- Two layer neural network: $f = W_2 \max(0, W_1 x)$
 - $x \in \mathcal{R}^D$
 - $W_1 \in \mathcal{R}^{H \times D}$
 - $W_2 \in \mathcal{R}^{C \times H}$
- In practice we will usually add a learnable bias at each layer as well
- Neural Network is a very broad term; these are more accurately called "fully-connected networks" or sometimes "multi-layer perceptrons" (MLP)

Neural networks: the original linear classifier

- Linear score function: $\mathbf{f} = \mathbf{W}\mathbf{x}$
 - $x \in \mathcal{R}^D$
 - $W \in \mathcal{R}^{C \times D}$
- Two layer neural network: $f = W_2 \max(0, W_1 x)$
 - $x \in \mathcal{R}^D$
 - $W_1 \in \mathcal{R}^{H \times D}$
 - $W_2 \in \mathcal{R}^{C \times H}$
- In practice we will usually add a learnable bias at each layer as well
- Neural Network is a very broad term; these are more accurately called "fully-connected networks" or sometimes "multi-layer perceptrons" (MLP)

Neural networks: the original linear classifier

- Linear score function: $\mathbf{f} = \mathbf{W}\mathbf{x}$
 - $x \in \mathcal{R}^D$
 - $W \in \mathcal{R}^{C \times D}$
- Two layer neural network: $f = W_2 \max(0, W_1 x)$
 - $x \in \mathcal{R}^D$
 - $W_1 \in \mathcal{R}^{H \times D}$
 - $W_2 \in \mathcal{R}^{C \times H}$
- In practice we will usually add a learnable bias at each layer as well
- Neural Network is a very broad term; these are more accurately called "fully-connected networks" or sometimes "multi-layer perceptrons" (MLP)

Neural networks: the original linear classifier

- Linear score function: $\mathbf{f} = \mathbf{W}\mathbf{x}$
 - $x \in \mathcal{R}^D$
 - $W \in \mathcal{R}^{C \times D}$
- Two layer neural network: $f = W_2 \max(0, W_1 x)$
 - $x \in \mathcal{R}^D$
 - $W_1 \in \mathcal{R}^{H \times D}$
 - $W_2 \in \mathcal{R}^{C \times H}$
- In practice we will usually add a learnable bias at each layer as well
- Neural Network is a very broad term; these are more accurately called "fully-connected networks" or sometimes "multi-layer perceptrons" (MLP)

Neural networks: why is max operator important?

- The function is called the activation function.
- Q: What if we try to build a neural network without one?
 - $f = W_2 W_1 x$
 - $W_3 = W_2 W_1 \in \mathcal{R}^{C \times D}$
 - $f = W_3 x$
- A: We end up with a linear classifier again!

Neural networks: why is max operator important?

- The function is called the activation function.
- Q: What if we try to build a neural network without one?
 - $f = W_2 W_1 x$
 - $W_3 = W_2 W_1 \in \mathcal{R}^{C \times D}$
 - $f = W_3 x$
- A: We end up with a linear classifier again!

Neural networks: why is max operator important?

- The function is called the activation function.
- Q: What if we try to build a neural network without one?
 - $f = W_2 W_1 x$
 - $W_3 = W_2 W_1 \in \mathcal{R}^{C \times D}$
 - $f = W_3 x$
- A: We end up with a linear classifier again!

Neural networks: why is max operator important?

- The function is called the activation function.
- Q: What if we try to build a neural network without one?
 - $f = W_2 W_1 x$
 - $W_3 = W_2 W_1 \in \mathcal{R}^{C \times D}$
 - $f = W_3 x$
- A: We end up with a linear classifier again!

Neural networks: why is max operator important?

- The function is called the activation function.
- Q: What if we try to build a neural network without one?
 - $f = W_2 W_1 x$
 - $W_3 = W_2 W_1 \in \mathcal{R}^{C \times D}$
 - $f = W_3 x$
- A: We end up with a linear classifier again!

Neural networks: why is max operator important?

- The function is called the activation function.
- Q: What if we try to build a neural network without one?
 - $f = W_2 W_1 x$
 - $W_3 = W_2 W_1 \in \mathcal{R}^{C \times D}$
 - $f = W_3 x$
- A: We end up with a linear classifier again!

Problem: How to compute gradients?

- $s = f(x; W_1, W_2) = W_2 \max(0, W_1 x)$ nonlinear score function
- $\mathcal{L}_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$ SVM loss
- $R(W) = \sum_k W_k^2$ regularization
- $L = \frac{1}{N} \sum_i \mathcal{L}_i + \lambda_1 R(W_1) + \lambda_2 R(W_2)$ Total loss
- If we can compute $\frac{\partial \mathcal{L}}{\partial W_1}$ and $\frac{\partial \mathcal{L}}{\partial W_2}$ then we can learn W_1 and W_2 .

Problem: How to compute gradients?

- $s = f(x; W_1, W_2) = W_2 \max(0, W_1 x)$ nonlinear score function
- $\mathcal{L}_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$ SVM loss
- $R(W) = \sum_k W_k^2$ regularization
- $L = \frac{1}{N} \sum_i \mathcal{L}_i + \lambda_1 R(W_1) + \lambda_2 R(W_2)$ Total loss
- If we can compute $\frac{\partial \mathcal{L}}{\partial W_1}$ and $\frac{\partial \mathcal{L}}{\partial W_2}$ then we can learn W_1 and W_2 .

Problem: How to compute gradients?

- $s = f(x; W_1, W_2) = W_2 \max(0, W_1 x)$ nonlinear score function
- $\mathcal{L}_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$ SVM loss
- $R(W) = \sum_k W_k^2$ regularization
- $L = \frac{1}{N} \sum_i \mathcal{L}_i + \lambda_1 R(W_1) + \lambda_2 R(W_2)$ Total loss
- If we can compute $\frac{\partial \mathcal{L}}{\partial W_1}$ and $\frac{\partial \mathcal{L}}{\partial W_2}$ then we can learn W_1 and W_2 .

Problem: How to compute gradients?

- $s = f(x; W_1, W_2) = W_2 \max(0, W_1 x)$ nonlinear score function
- $\mathcal{L}_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$ SVM loss
- $R(W) = \sum_k W_k^2$ regularization
- $L = \frac{1}{N} \sum_i^N \mathcal{L}_i + \lambda_1 R(W_1) + \lambda_2 R(W_2)$ Total loss
- If we can compute $\frac{\partial \mathcal{L}}{\partial W_1}$ and $\frac{\partial \mathcal{L}}{\partial W_2}$ then we can learn W_1 and W_2 .

Problem: How to compute gradients?

- $s = f(x; W_1, W_2) = W_2 \max(0, W_1 x)$ nonlinear score function
- $\mathcal{L}_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$ SVM loss
- $R(W) = \sum_k W_k^2$ regularization
- $L = \frac{1}{N} \sum_i^N \mathcal{L}_i + \lambda_1 R(W_1) + \lambda_2 R(W_2)$ Total loss
- If we can compute $\frac{\partial \mathcal{L}}{\partial W_1}$ and $\frac{\partial \mathcal{L}}{\partial W_2}$ then we can learn W_1 and W_2 .

Problem: How to compute gradients?

- $s = f(x; W_1, W_2) = W_2 \max(0, W_1 x)$ nonlinear score function
- $\mathcal{L}_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$ SVM loss
- $R(W) = \sum_k W_k^2$ regularization
- $L = \frac{1}{N} \sum_i^N \mathcal{L}_i + \lambda_1 R(W_1) + \lambda_2 R(W_2)$ Total loss
- If we can compute $\frac{\partial \mathcal{L}}{\partial W_1}$ and $\frac{\partial \mathcal{L}}{\partial W_2}$ then we can learn W_1 and W_2 .

Problem: How to compute gradients?

- $s = f(x; W_1, W_2) = W_2 \max(0, W_1 x)$ nonlinear score function
- $\mathcal{L}_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$ SVM loss
- $R(W) = \sum_k W_k^2$ regularization
- $L = \frac{1}{N} \sum_i^N \mathcal{L}_i + \lambda_1 R(W_1) + \lambda_2 R(W_2)$ Total loss
- If we can compute $\frac{\partial \mathcal{L}}{\partial W_1}$ and $\frac{\partial \mathcal{L}}{\partial W_2}$ then we can learn W_1 and W_2 .

Problem: How to compute gradients?

- $s = f(x; W_1, W_2) = W_2 \max(0, W_1 x)$ nonlinear score function
- $\mathcal{L}_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$ SVM loss
- $R(W) = \sum_k W_k^2$ regularization
- $L = \frac{1}{N} \sum_i^N \mathcal{L}_i + \lambda_1 R(W_1) + \lambda_2 R(W_2)$ Total loss
- If we can compute $\frac{\partial \mathcal{L}}{\partial W_1}$ and $\frac{\partial \mathcal{L}}{\partial W_2}$ then we can learn W_1 and W_2 .

Problem: How to compute gradients?

- $s = f(x; W_1, W_2) = W_2 \max(0, W_1 x)$ nonlinear score function
- $\mathcal{L}_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$ SVM loss
- $R(W) = \sum_k W_k^2$ regularization
- $L = \frac{1}{N} \sum_i^N \mathcal{L}_i + \lambda_1 R(W_1) + \lambda_2 R(W_2)$ Total loss
- If we can compute $\frac{\partial \mathcal{L}}{\partial W_1}$ and $\frac{\partial \mathcal{L}}{\partial W_2}$ then we can learn W_1 and W_2 .

Stochastic gradient descent

$$\min_{\mathbf{w}} \sum_i L(NN(\mathbf{x}_i, \mathbf{w}), y_i)$$

Stochastic gradient descent

Given a training set $S = \{(\mathbf{x}_i, y_i)\}, \mathbf{x} \in \mathbb{R}^d$

1. Initialize parameters \mathbf{w}
2. For epoch = 1 ... T:
 1. Shuffle the training set
 2. For each training example $(\mathbf{x}_i, y_i) \in S$:
 - Treat this example as the entire dataset
 - Compute the gradient of the loss $\nabla L(NN(\mathbf{x}_i, \mathbf{w}), y_i)$ using backpropagation
 - Update: $\mathbf{w} \leftarrow \mathbf{w} - \gamma_t \nabla L(NN(\mathbf{x}_i, \mathbf{w}), y_i)$

The objective is **not convex**.
Initialization can be important

3. Return \mathbf{w}

γ_t : learning rate,
many tweaks possible

Lets Calculate

(Bad) Idea: Derive $\nabla_W L$ on paper

$$s = f(x; W) = Wx$$

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

$$= \sum_{j \neq y_i} \max(0, W_{j,:} \cdot x + W_{y_i,:} \cdot x + 1)$$

$$L = \frac{1}{N} \sum_{i=1}^N L_i + \lambda \sum_k W_k^2$$

$$= \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, W_{j,:} \cdot x + W_{y_i,:} \cdot x + 1) + \lambda \sum_k W_k^2$$

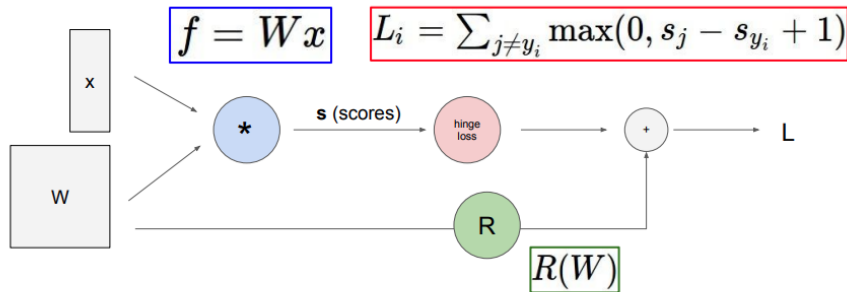
$$\nabla_W L = \nabla_W \left(\frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, W_{j,:} \cdot x + W_{y_i,:} \cdot x + 1) + \lambda \sum_k W_k^2 \right)$$

Problem: Very tedious: Lots of matrix calculus, need lots of paper

Problem: What if we want to change loss? E.g. use softmax instead of SVM? Need to re-derive from scratch = (

Problem: Not feasible for very complex models!

Better Idea: Computational graphs + Backpropagation



Solution: Backpropagation!!!!

Examples

Backpropagation: a simple example

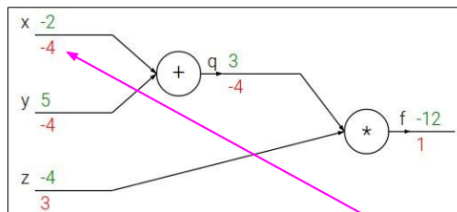
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



Chain rule:

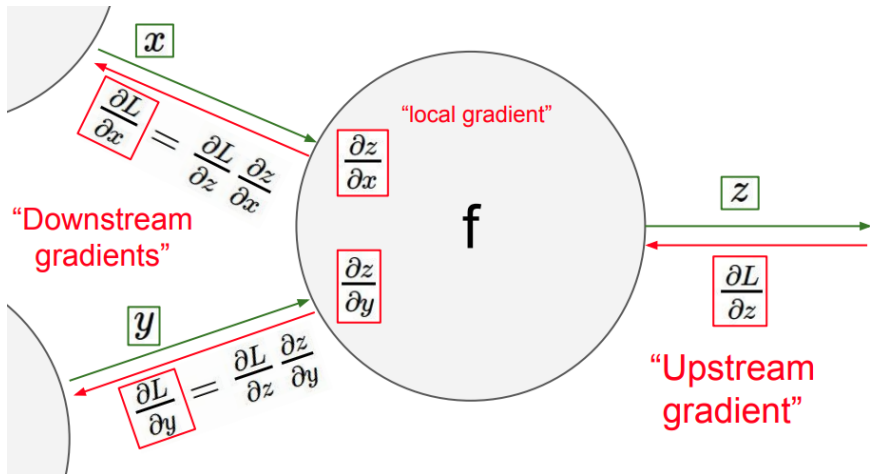
$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

Upstream
gradient

Local
gradient

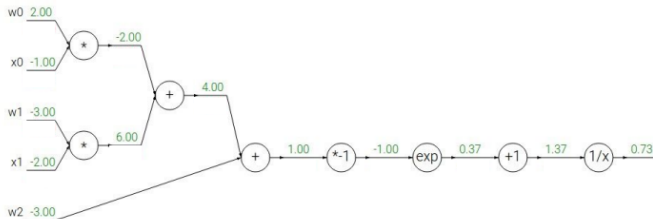
$$\frac{\partial f}{\partial x}$$

Examples



Examples

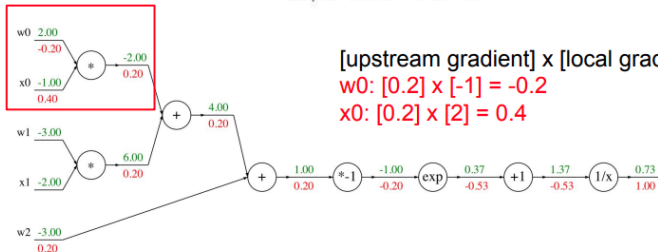
Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$



$f(x) = e^x$	\rightarrow	$\frac{df}{dx} = e^x$		$f(x) = \frac{1}{x}$	\rightarrow	$\frac{df}{dx} = -1/x^2$
$f_a(x) = ax$	\rightarrow	$\frac{df}{dx} = a$		$f_c(x) = c + x$	\rightarrow	$\frac{df}{dx} = 1$

Examples

Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2x_2)}}$



[upstream gradient] x [local gradient]

$w_0: [0.2] \times [-1] = -0.2$

$x_0: [0.2] \times [2] = 0.4$

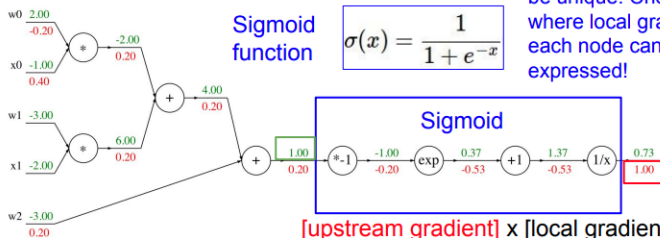
$f(x) = e^x$	\rightarrow	$\frac{df}{dx} = e^x$		$f(x) = \frac{1}{x}$	\rightarrow	$\frac{df}{dx} = -1/x^2$
$f_a(x) = ax$	\rightarrow	$\frac{df}{dx} = a$		$f_c(x) = c + x$	\rightarrow	$\frac{df}{dx} = 1$

Examples

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

Computational graph representation may not be unique. Choose one where local gradients at each node can be easily expressed!



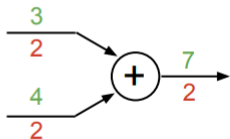
[upstream gradient] x [local gradient]
 $[1.00] \times [(1 - 1/(1+e^1)) (1/(1+e^1))] = 0.2$

Sigmoid local gradient:

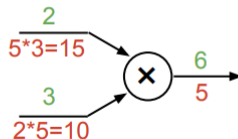
$$\frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = \left(\frac{1 + e^{-x} - 1}{1 + e^{-x}} \right) \left(\frac{1}{1 + e^{-x}} \right) = (1 - \sigma(x)) \sigma(x)$$

Patterns in gradient flow

add gate: gradient distributor



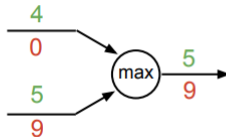
mul gate: “swap multiplier”



copy gate: gradient adder

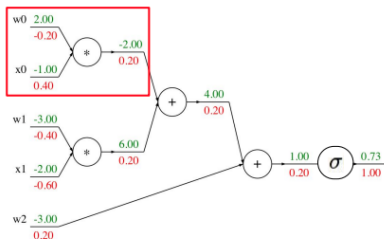


max gate: gradient router



Backward Code

Backprop Implementation: “Flat” code



Forward pass:
Compute output

```
def f(w0, x0, w1, x1, w2):
```

```
    s0 = w0 * x0
```

```
    s1 = w1 * x1
```

```
    s2 = s0 + s1
```

```
    s3 = s2 + w2
```

```
    L = sigmoid(s3)
```

```
grad_L = 1.0
```

```
grad_s3 = grad_L * (1 - L) * L
```

```
grad_w2 = grad_s3
```

```
grad_s2 = grad_s3
```

```
grad_s0 = grad_s2
```

```
grad_s1 = grad_s2
```

```
grad_w1 = grad_s1 * x1
```

```
grad_x1 = grad_s1 * w1
```

```
grad_w0 = grad_s0 * x0
```

```
grad_x0 = grad_s0 * w0
```

Multiply gate

Vector derivatives

Scalar to Scalar

$$x \in \mathbb{R}, y \in \mathbb{R}$$

Regular derivative:

$$\frac{\partial y}{\partial x} \in \mathbb{R}$$

If x changes by a small amount, how much will y change?

Vector to Scalar

$$x \in \mathbb{R}^N, y \in \mathbb{R}$$

Derivative is **Gradient**:

$$\frac{\partial y}{\partial x} \in \mathbb{R}^N \quad \left(\frac{\partial y}{\partial x} \right)_n = \frac{\partial y}{\partial x_n}$$

For each element of x , if it changes by a small amount then how much will y change?

Vector to Vector

$$x \in \mathbb{R}^N, y \in \mathbb{R}^M$$

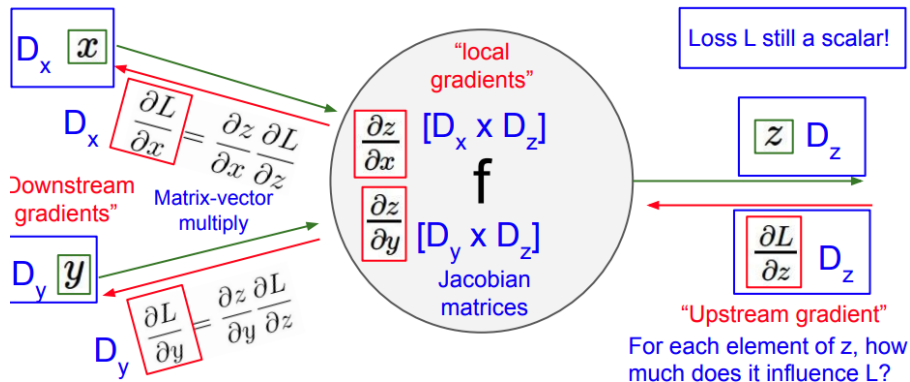
Derivative is **Jacobian**:

$$\frac{\partial y}{\partial x} \in \mathbb{R}^{N \times M} \quad \left(\frac{\partial y}{\partial x} \right)_{n,m} = \frac{\partial y_m}{\partial x_n}$$

For each element of x , if it changes by a small amount then how much will each element of y change?

Vector derivatives

Backprop with Vectors



Vector derivatives: Example

Jacobian is **sparse**:
off-diagonal entries
always zero! Never
explicitly form
Jacobian -- instead
use **implicit**
multiplication

4D input x:

$\begin{bmatrix} 1 \\ -2 \\ 3 \\ -1 \end{bmatrix}$

$f(x) = \max(0, x)$
(*elementwise*)

4D output z:

$\begin{bmatrix} 1 \\ 0 \\ 3 \\ 0 \end{bmatrix}$

4D dL/dx:

$\begin{bmatrix} 4 \\ 0 \\ 5 \\ 0 \end{bmatrix}$

$[dz/dx] [dL/dz]$

$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 4 \\ -1 \\ 5 \\ 9 \end{bmatrix}$

4D dL/dz:

$\begin{bmatrix} 4 \\ -1 \\ 5 \\ 9 \end{bmatrix}$

Upstream
gradient

Vector derivatives: Example

Jacobian is **sparse**:
off-diagonal entries
always zero! Never
explicitly form
Jacobian -- instead
use **implicit**
multiplication

4D input x:

$\begin{bmatrix} 1 \\ -2 \\ 3 \\ -1 \end{bmatrix}$

$f(x) = \max(0, x)$
(*elementwise*)

4D output z:

$\begin{bmatrix} 1 \\ 0 \\ 3 \\ 0 \end{bmatrix}$

4D dL/dx :

$\begin{bmatrix} 4 \\ 0 \\ 5 \\ 0 \end{bmatrix}$

$[dz/dx] [dL/dz]$

$$\left(\frac{\partial L}{\partial x}\right)_i = \begin{cases} \left(\frac{\partial L}{\partial z}\right)_i & \text{if } x_i > 0 \\ 0 & \text{otherwise} \end{cases}$$

4D dL/dz :

$\begin{bmatrix} 4 \\ -1 \\ 5 \\ 9 \end{bmatrix}$

Upstream
gradient

Metrics derivatives: Example

Jacobian is **sparse**:
off-diagonal entries
always zero! Never
explicitly form
Jacobian -- instead
use **implicit**
multiplication

4D input x:

$\begin{bmatrix} 1 \\ -2 \\ 3 \\ -1 \end{bmatrix}$

$f(x) = \max(0, x)$
(*elementwise*)

4D output z:

$\begin{bmatrix} 1 \\ 0 \\ 3 \\ 0 \end{bmatrix}$

4D dL/dx:

$\begin{bmatrix} 4 \\ 0 \\ 5 \\ 0 \end{bmatrix}$

$[dz/dx] [dL/dz]$

$$\left(\frac{\partial L}{\partial x}\right)_i = \begin{cases} \left(\frac{\partial L}{\partial z}\right)_i & \text{if } x_i > 0 \\ 0 & \text{otherwise} \end{cases}$$

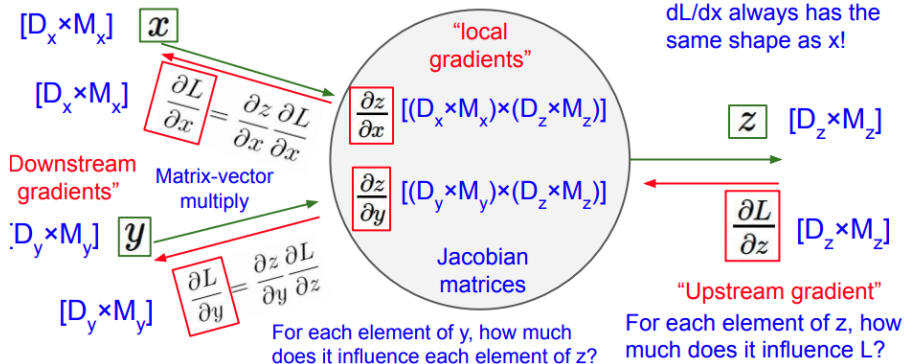
4D dL/dz:

$\begin{bmatrix} 4 \\ -1 \\ 5 \\ 9 \end{bmatrix}$

Upstream
gradient

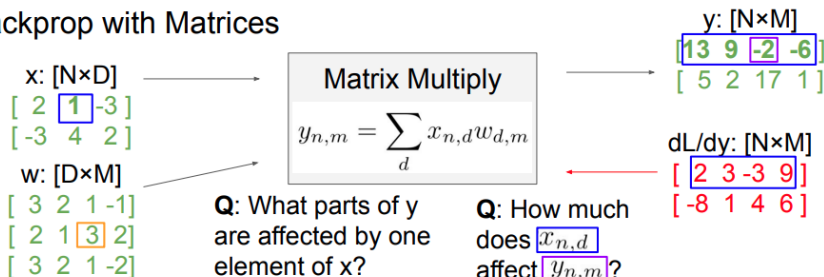
Metrics derivatives: Example

Backprop with Matrices (or Tensors)



Metrics derivatives: Example

Backprop with Matrices



Q: What parts of y are affected by one element of x ?

A: $x_{n,d}$ affects the whole row $y_{n,\cdot}$.

Q: How much does $x_{n,d}$ affect $y_{n,m}$?

A: $w_{d,m}$

$$\frac{\partial L}{\partial x_{n,d}} = \sum_m \frac{\partial L}{\partial y_{n,m}} \frac{\partial y_{n,m}}{\partial x_{n,d}} = \sum_m \frac{\partial L}{\partial y_{n,m}} w_{d,m}$$

Metrics derivatives: Example

Backprop with Matrices

$$x: [N \times D]$$

$$\begin{bmatrix} 2 & 1 & -3 \\ -3 & 4 & 2 \end{bmatrix}$$

$$w: [D \times M]$$

$$\begin{bmatrix} 3 & 2 & 1 & -1 \\ 2 & 1 & 3 & 2 \\ 3 & 2 & 1 & -2 \end{bmatrix}$$

Matrix Multiply

$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

$$y: [N \times M]$$

$$\begin{bmatrix} 13 & 9 & -2 & -6 \\ 5 & 2 & 17 & 1 \end{bmatrix}$$

$$dL/dy: [N \times M]$$

$$\begin{bmatrix} 2 & 3 & -3 & 9 \\ -8 & 1 & 4 & 6 \end{bmatrix}$$

By similar logic:

$$[N \times D] \quad [N \times M] \quad [M \times D]$$

$$\frac{\partial L}{\partial x} = \left(\frac{\partial L}{\partial y} \right) w^T$$

$$[D \times M] \quad [D \times N] \quad [N \times M]$$

$$\frac{\partial L}{\partial w} = x^T \left(\frac{\partial L}{\partial y} \right)$$

These formulas are easy to remember: they are the only way to make shapes match up!

Optimization and activation

Next Class..